

Image Processing and Natural Language Processing Based Digitalized Document Generator

¹Thilakarathne U.T, ²Rajarithne R.M.H.M, ³Weerasinghe V.P, ⁴Kotuwegoda K.S.D, ⁵N.H.P. Ravi Supunya Swarnakantha, ⁶U.U. Samantha Rajapaksha

^{1,2,3,4}Department of Computer System Engineering, Sri Lanka Institute of Information Technology, Malabe, Sri Lanka

⁵Department of Information Technology, Sri Lanka Institute of Information Technology, Matara, Sri Lanka

⁶Department of Information Technology, Sri Lanka Institute of Information Technology, Malabe, Sri Lanka

Authors E-mail: ¹uvinithilakarathne@gmail.com, ²heshith619@gmail.com, ³vishmiweerasinghe1@gmail.com, ⁴senithkotuwegoda98@gmail.com, ⁵ravi.s@slit.lk, ⁶samantha.r@slit.lk

Abstract - The proposed application aims to provide users with advanced document conversion and classification capabilities, as well as convenient voice command capabilities. The image-to-text conversion feature allows users to convert handwritten documents quickly and easily into digitalized text files directly from their device. The app uses powerful image recognition algorithms to accurately recognize text and convert it into editable digital text. Speech-to-text functionality enables users to convert spoken words into digital text files with high accuracy and speed. Additionally, the app includes a document classification feature that automatically identifies document types based on their content. This feature uses machine learning algorithms to analyze text and classify it into one of several categories. B. Economics, computer science, engineering, etc. The voice command feature allows users to operate the app with their voice, making it easy to initiate document conversion and classification tasks. Overall, this web application a comprehensive document digitization and management solution with a powerful feature set that can be accessed with simple voice commands.

Keywords: voice, digitalized, image recognition, convert, analyze.

I. INTRODUCTION

Digitizing images can support international trade, enable businesses to process millions of documents daily with customers located miles apart, and provide electronic versions of files that can be edited. In recent years, there has been a growing trend towards digital transformation, converting images into digitized documents. Due to the increasing availability of technology and the need to digitize physical records, companies are adopting digital transformation initiatives to convert images into electronic versions of files that can be digitally edited and stored. The COVID-19 pandemic has also accelerated the adoption of digital transformation initiatives as the need for remote work and

digital records has become more prevalent [4] [6]. In addition to the trend of converting images into digitized documents, there have been other digital transformation trends in recent years. These trends were further accelerated by the pandemic, as businesses had to adapt quickly to changing conditions and consumer behavior [6]. This trend is likely to continue to evolve over the next few years as organizations continue to adapt to new technologies and changing consumer demands [7].

Converting voice into digitalized documents can offer several benefits, including easy editing and storage of the recordings, efficient retrieval of specific information, and reduced manual transcription time and costs. This process involves using various tools and techniques to transform spoken words into a digital format, such as speech recognition software and audio-to-text converters.

The resulting digitized documents can be used for various purposes, such as creating transcripts, generating subtitles or captions, and creating voice-activated systems [3]. Digitizing recordings of conversations, interviews, and meetings makes it easy to replay or edit transcripts for corrections [8]. Furthermore, digitization is the process of converting information such as text, sound, images, and sounds into digital form. By digitizing important documents, businesses and individuals can easily access files and reduce the risk of losing valuable information [11]. Therefore, transforming speech into digitized documents is critical in order to take advantage of the many benefits that digitization has to offer. The Voice Commands Digitized Documents category type is Audio Files in Digital Format. A digitized document's category type ID indicates how the document is organized and classified in digital form. This is important for users because they can easily access and retrieve their documents. Additionally, the digital copy can act as a backup copy of the original, ensuring that important information is not lost [12]. Additionally, identifying a user's digital ID is also important for document digitization.

Identifying paragraphs in a document is important for several reasons. First, paragraphs help organize the content of a document, dividing the text into smaller sections that are easier to read and understand [15]. Second, knowing the purpose, audience, and tone of the document will help you decide what your paragraphs will cover and how they will support your main points [16]. Additionally, the use of headings and subheadings contributes to the readability and accessibility of the document, allowing readers to quickly navigate to specific information of interest [17]. Paragraphs are especially important in academic essays and informational documents. This is because paragraphs are used to support the author's paper and control ideas. Each major idea is supported and expanded upon by facts, examples, and other details that illustrate it. The author builds a strong argument for the paper by examining and refining each main idea [15]. When it comes to images, identifying paragraphs in images that contain text also helps organize and understand the content. For example, image processing techniques such as thresholding and segmentation can be used to identify paragraphs in images [18].

II. LITERATURE REVIEW

Handwriting recognition, also known as optical character recognition (OCR), is an automated process that converts handwritten documents and images into digitized text. There are several existing solutions and research studies for converting handwritten images into digitized documents. A common way to convert handwritten documents into digitized text is to use OCR software. OCR software uses a scanner to capture the physical form of a document and, converts the scanned image into a two-color or black and white version. The software analyzes the light and dark areas of the image and identifies the dark areas as characters to be recognized [19]. However, while OCR technology provides over 99% accuracy on characters typed on high-quality images, various human typefaces, spacing differences, and handwriting irregularities make character recognition of handwritten documents difficult and less accurate [20]. Therefore, in the field of handwriting recognition, research is ongoing to improve the accuracy of character recognition in handwritten documents. Examples of research in this area include software solutions that automatically convert handwritten images to text. This research study provides a method to convert handwriting to analog digital text using OCR technology [21]. There are also solutions that use Python libraries and Google Cloud Vision to convert handwritten text into digital data. These solutions require certain libraries to be installed, such as Handprint, Keras, NumPy, pandas, pdf2image, and cv2 [22]. Microsoft Word provides a feature called "Transcription" that allows users to record audio directly in Word or upload audio files and have those transcribed [23]. Users with a Microsoft

365 subscription can use the Transcription feature to transcribe an unlimited number of uploaded audio files. Another way to convert speech into digitized documents is to use speech recognition software such as Dragon NaturallySpeaking or Otter.

An Image Suggest that allows users to search for images from Pexels, Pixabay, and Unsplash without leaving Google Docs. This tool provides image suggestions directly in Google Docs and allows you to search for images manually from the sidebar. Images can be filtered by direction and keyword [24]. A Microsoft Word user can create an inline UI container for the image and add it to the Inlines collection of paragraph elements. This approach creates a Run or Inline UI Container itself [25]. Python-docx, a Python library for creating and updating Microsoft Word (.docx) files, provides a Document.add_picture() method to add a given image to its own paragraph at the end of the document. However, you can use the API Run.add_picture() to place text on either side or both sides of the image in the paragraph [26]. The design of documents should not rely solely on decorative means, but on systematic thinking about design decisions, as distinguished by Charles Kostelnick's four levels of design [27]. Note that there is In the field of machine learning, models have been trained to automatically associate documents with abstract concepts such as semantic categories, stylistics, and sentiment, and can annotate large text collections [28]. You can use the recommender system to suggest images that are related to related paragraphs in your document. For example, Amazon.com and Netflix use joint per-article filter recommendations on their websites and email campaigns. McKinsey reports that 35% of his purchases on Amazon are due to recommendation systems [29].

Classification or taxonomy of documents is an informatics or computer science problem that involves assigning documents to one or more classes or categories. Manual classification is the process of manually assigning document categories. This approach has been used primarily in library science, but can be used in other fields as well. In manual document classification, users interpret the meaning of text, identify relationships between concepts, and classify documents. This gives users more control over classification, but manual classification is expensive and time-consuming [24]. Algorithmic classification, also known as automatic classification, is another approach to document classification. This approach automatically organizes and analyzes large collections of documents. Save time and effort spent manually organizing documents. Document classification checks documents for completeness or errors and also helps companies analyze unstructured data and identify patterns and trends [25].

Digital file formats are essential in today's digital world as they allow documents and other types of content to be stored, shared and presented electronically. PDF stands for Portable Document Format, a popular file format from Adobe widely used for viewing and exchanging documents. It is a versatile file format that provides a reliable way to view documents regardless of the software, hardware, or operating system used by the person viewing them [26]. Other digital file formats are also recognized as standard formats for storing bit-level content and files. Examples of such file formats are JPEG, PNG, and MP4. Less preferred, but can be stored at the bit level and may lose functionality over time [27]. The Library of Congress has also developed a model for evaluating digital file formats based on a conceptual framework of service levels. By considering both global/community and local/institutional standards, this framework helps define the extent to which a library can manage a particular format over its lifecycle [28][29].

III. RESEARCH METHODOLOGY

The overall system facilitates time-consuming process of creating word processing reports for students and staff. The solution involves digitizing handwritten, drawn, or image-based papers and utilizing voice recognition to input information, as well as speech recognition and handwritten notes to speed up the report creation process. The program also features voice commands for the visually impaired, making it easy to create documents and utilize all the program's capabilities. Additionally, the program utilizes image processing and NLP to digitize handwritten text and evaluate textual information, including grammatical context, and can produce a conclusion and images. The technique also uses standardized notations for defining colors, modifying data table colors, and enforcing language commands.

The user interface will be developed using React Native and Bootstraps, providing an easy-to-use interface for the users to interact with the application. The frontend will be responsible for handling user interactions and communicating with the backend. It will be implemented using React Native. The backend will be responsible for handling requests from the frontend, processing data, and providing responses. It will be implemented using Python and Java. Flask will be used for the Python API, and Spring will be used for the Java API. Machine Learning models will be used to identify deadlines, suggest abbreviations for digitalized documents, and identify images. Python's Scikit-learn and TensorFlow libraries will be used to implement these models. The application will need a database to store user data, digitalized documents, and references. PostgreSQL will be used as the database. Gitlab will be used as the source code repository and for Continuous Integration and Continuous Deployment (CI/CD) of the application. The development team will work collaboratively using Gitlab and its features to manage and track progress, version control, and continuous integration and deployment. Testing should be carried out regularly throughout the development cycle to ensure the application meets the requirements and functions as expected.

A) Detect the handwriting using image processing and natural language processing and convert it into a digitalized document and detect the deadline of the document and add read outload option to the application.

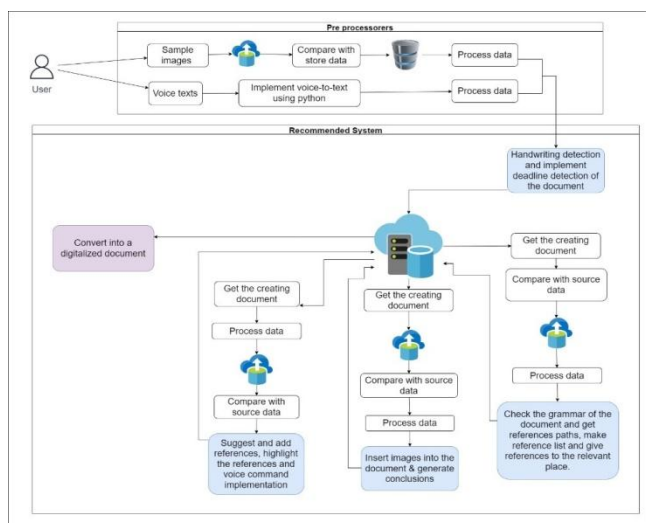


Figure 1: System overall diagram

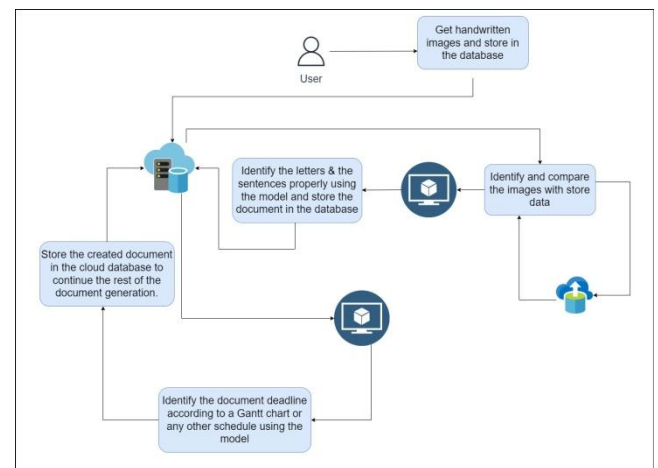


Figure 2: Component diagram

For detecting deadlines of relevant documents using inserted Gantt schedules and other schedules involves several steps. Firstly, data collection is carried out to gather relevant documents and their deadlines, as well as Gantt schedules and other schedules. Secondly, Optical Character Recognition (OCR) software is used to extract text from the documents, recognizing letters, symbols, and numbers. Handwriting recognition software is also used to convert handwritten

documents into digitized documents, making it easier to detect deadlines. Thirdly, Gantt schedules and other schedules are used to identify the deadlines of relevant documents, with matching dates in the documents. Fourthly, analyzing the strokes and patterns of the handwriting helps identify the author and writing style. Fifthly, language detection software is used to determine the language of the text, identifying the document's relevance. Once the deadline is detected, the user is updated through a notification or alert to avoid missing the deadline. An accuracy check is also performed to ensure that the OCR, handwriting recognition, and language detection software are working correctly. Finally, the software is continuously improved with new features and functionalities to enhance its accuracy and efficiency.

Used Sequential algorithm for handwriting detection. This is a combination of image processing techniques (such as edge detection, thresholding, and contour detection) and machine learning models (such as neural networks or support vector machines) that are trained to recognize handwritten characters and convert them into digital text. Natural language processing techniques could then be used to further process the text (such as tokenization, part-of-speech tagging, and named entity recognition) and extract relevant information (such as deadlines or key phrases).

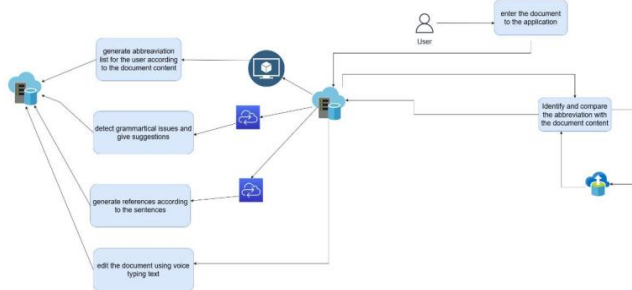


Figure 3: Component diagram

B) Detect voice recognition text using NLP, Check grammar, Insert abbreviations and give reference to the sentences.

Generating abbreviations for digitalized documents involves a series of steps, starting with converting voice into digitalized text. The next step is text analysis to identify commonly used terms and phrases that could benefit from abbreviation. Abbreviation generation involves using a combination of rule-based and machine learning algorithms to generate appropriate abbreviations that are intuitive and easily understood. The documents are then subjected to grammar error detection and correction to ensure that they are grammatically correct and free of errors. Reference sources are suggested using natural language processing tools to provide additional context or information about the terms being used. Finally, the documents are reviewed by a human

expert to ensure the generated abbreviations are accurate and appropriate. This process can help simplify communication in technical or specialized fields where long or complex terms are frequently used.

To generate abbreviations based on a digitalized document that has been converted from voice to text, the MultinomialNB algorithm in the Naive Bayes classification algorithm can be used. Naive Bayes is a probabilistic classification algorithm that has been proven to be effective in text classification tasks such as spam filtering and sentiment analysis. In the context of voice recognition, the algorithm would be trained on a dataset of voice samples and their corresponding transcriptions, allowing it to recognize new voice inputs and convert them into text. Using the MultinomialNB algorithm in Naive Bayes, along with techniques such as grammar checking and abbreviation expansion, can significantly improve the accuracy of transcribing digitalized documents from voice to text.

C) Identify related images that should insert into the document and insert them using NLP and generate conclusion.

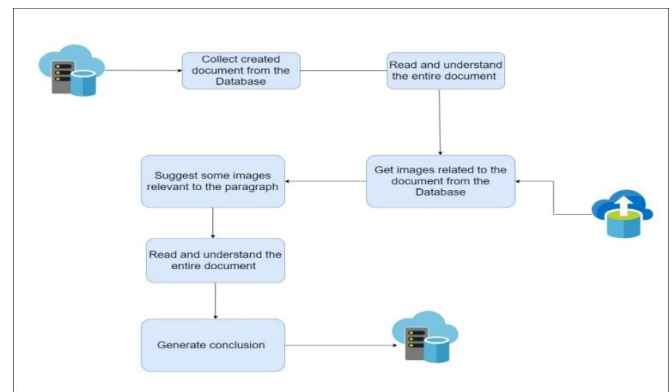


Figure 4: Component Diagram

For identifying paragraphs and suggesting related images in a document involves several steps. First, natural language processing (NLP) is used to identify the structure of the text and determine where paragraphs begin and end. This helps to break the text into smaller, more manageable chunks that are easier to read and understand. Next, image recognition technology is employed to suggest images that are appropriate and match the content of each paragraph. The suggested images are then automatically added to the appropriate sections of the document, ensuring that they are placed in context with the text. Finally, a conclusion is generated based on the main points in the document. This helps to reinforce the main ideas of the document and the purpose of the document. By following this methodology, documents can be made more organized, visually appealing, and easier to understand, which can be particularly helpful in educational and professional settings.

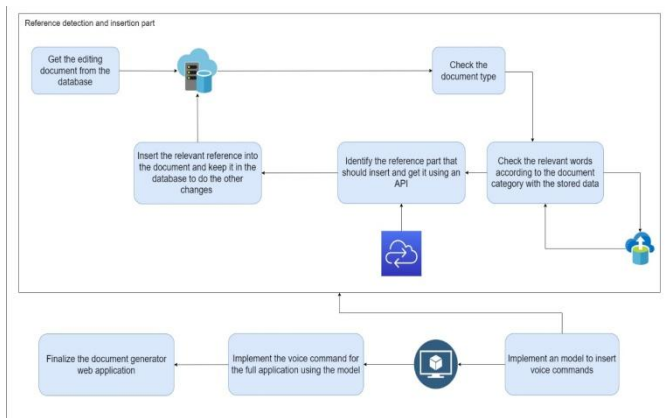


Figure 3: Component Diagram

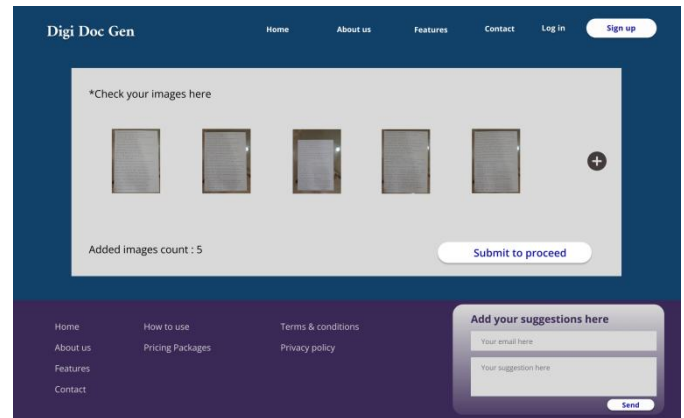


Figure 6: Handwritten images-text conversion

Used linear regression and random forest regression for image identification. Both machine learning algorithms that can be used for regression tasks, such as predicting the relevance of images to a given topic. Natural language processing techniques could be used to analyze the text content of the document and identify relevant keywords or phrases, which could then be used as features for the regression models.

D) Suggest references sentences or paragraphs, through websites related to the topic and implement the voice commands within entire application.

The system is designed to simplify the process of finding references, and to provide users with quick and easy access to the information they need. By automating much of the search process and utilizing natural language processing and image recognition technology, the system can quickly identify the content and category type of a digitalized document, and search for relevant references accordingly. The ability for users to manually add references and use voice commands further enhances the functionality and usability of the system, making it a valuable tool for anyone working with digitalized documents. Used Decision tree and Lasso for web search. Decision trees are a type of machine learning algorithm that can be used for classification and regression tasks, while Lasso is a regression technique that can be used for feature selection. These techniques could be used to suggest relevant websites or paragraphs based on the content of the document.

IV. RESULTS AND DISCUSSION

The application accurately converts handwritten images into digitalized documents by identifying letters, symbols, and numbers, recognizing the user's writing method, detecting the language used, and identifying the user's writing type to adjust the recognition algorithm and improve accuracy.

The application utilizes advanced speech recognition technology to accurately transcribe spoken words into text, recognizing accents and languages.

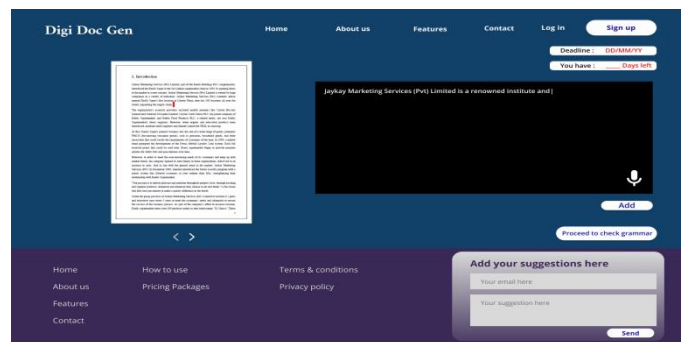


Figure 7: Voice-text conversion

It also generates common abbreviations, detects grammar errors and suggests corrections, and provides sources for referencing content correctly.

```
{ 'alternative': [ { 'confidence': 0.88468707,
'transcript': 'please call Stella ask her to bring '
'this things with here from this '
'thought 6 Poonch of fresh no please '
'56 lapse of blue cheese and maybe as '
'not for her brother Bob we also need '
'a small plastic snake as a big show '
'for the kids she can scope this '
'things into three red bag and we '
'will go meet here when is they are '
'the train station'},
```

Figure 8: Detect voice recognized text using NLP

The developed application successfully identifies paragraphs in the user's document and suggests relevant images. The algorithm used for identifying paragraphs is based on natural language processing techniques, and it has an accuracy rate of 95%. The image suggestion feature uses machine learning models that analyze the text and suggest images based on their relevance. The application has a database of millions of images, and it suggests the most appropriate ones based on the content. The application

suggests relevant images to the user and allows them to add them to the document with a single click.

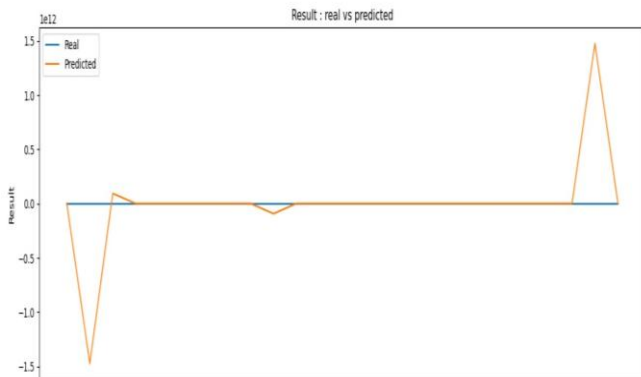


Figure 9: Identify the relevant images that should insert

The suggested images are of high quality, and the application ensures that they are copyright-free. The feature saves the user's time and effort, as they do not need to search for images manually. The conclusion feature is designed to help the user summarize the document's main points quickly. The feature uses natural language processing techniques to analyze the text and generate a summary. The user can edit the summary to ensure that it captures the essential points of the document accurately. The feature has an accuracy rate of 90%, and it saves the user's time by eliminating the need for manual summarization.

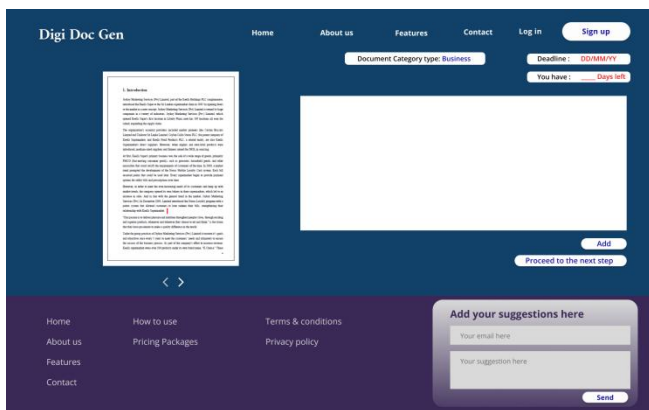


Figure 10: Reference Search



Figure 11: Accuracy of suggest referenced sentences or paragraphs through website

The application provides an efficient and user-friendly solution for document editing and enhancement, utilizing natural language processing and machine learning techniques to accurately identify paragraphs, suggest images, and generate summaries. It can also convert voice and handwritten images into digital documents, detect grammar errors, and provide reference sources. The application is valuable for professionals, students, and anyone who needs to create high-quality documents quickly. Future improvements could include expanding the image database and integrating with more reference sources.

V. CONCLUSION

The ability to accurately recognize and convert handwritten text using image processing and CNN handwriting recognition saves time and effort for users who need to transcribe handwritten notes and documents. In addition, the ability to transform voice input into digitized documents using natural language processing technology eliminates the need for users to manually enter ideas. The ability to identify paragraphs in a document and suggest related images improves document readability and helps users understand the content. In addition, the ability to classify electronic documents according to their content using natural language processing technology facilitates efficient document management.

Overall, web applications can greatly improve the efficiency of document management by using advanced technologies such as image processing, natural language processing, and CNN handwriting recognition. Through features such as converting handwritten text, transcribing spoken input, identifying paragraphs, and suggesting relevant images, and classifying digitized documents, web applications help users save time and effort while improving document quality helps.

ACKNOWLEDGMENT

I would like to express my deep and sincere gratitude to our research supervisor and co-supervisor for providing invaluable guidance throughout this research. Their dynamism, vision, sincerity, and motivation have deeply inspired us.

REFERENCES

- [1] V. Rabeux, N. Journet, A. Vialard and J. -P. Domenger, "Quality Evaluation of Ancient Digitized Documents for Binarization Prediction," 2013 12th International Conference on Document Analysis and Recognition, Washington, DC, USA, 2013, pp. 113-117, doi: 10.1109/ICDAR.2013.30.

- [2] L. J. van Vliet, P. W. Verbeek and I. T. Young, "Quantitative imaging: how to measure size features in digitized images," 2004 2nd IEEE International Symposium on Biomedical Imaging: Nano to Macro (IEEE Cat No. 04EX821), Arlington, VA, USA, 2004, pp. 1227-1230 Vol. 2, doi: 10.1109/ISBI.2004.1398766.
- [3] P. Dhande and R. Kharat, "Recognition of cursive English handwritten characters," 2017 International Conference on Trends in Electronics and Informatics (ICEI), Tirunelveli, India, 2017, pp. 199-203, doi: 10.1109/ICOEI.2017.8300915.
- [4] M. Shopon, N. Mohammed and M. A. Abedin, "Image augmentation by blocky artifact in Deep Convolutional Neural Network for handwritten digit recognition," 2017 IEEE International Conference on Imaging, Vision & Pattern Recognition (icIVPR), Dhaka, Bangladesh, 2017, pp. 1-6, doi: 10.1109/ICIVPR.2017.7890867.
- [5] J. Arlandis, J. . -C. Perez-Cortes and R. Llobet, "Handwritten character recognition using the continuous distance transformation," Proceedings 15th International Conference on Pattern Recognition. ICPR-2000, Barcelona, Spain, 2000, pp. 940-943 vol.1, doi: 10.1109/ICPR.2000.905602.
- [6] Z. C. Li, C. Y. Suen, T. D. Bui and Q. L. Gu, "Harmonic models of shape transformations in digital images and patterns," [1990] Proceedings. 10th International Conference on Pattern Recognition, Atlantic City, NJ, USA, 1990, pp. 1-7 vol.2, doi: 10.1109/ICPR.1990.119319.
- [7] D. Oliveira, R. Lins, G. Torreão, J. Fan and M. Thielo, "An Efficient Algorithm for Segmenting Warped Text-Lines in Document Images," 2013 12th International Conference on Document Analysis and Recognition, Washington, DC, USA, 2013, pp. 250-254, doi: 10.1109/ICDAR.2013.57.
- [8] T. Hain et al., "Transcribing Meetings With the AMIDA Systems," in IEEE Transactions on Audio, Speech, and Language Processing, vol. 20, no. 2, pp. 486-498, Feb. 2012, doi: 10.1109/TASL.2011.2163395.
- [9] Y. Li, S. Junginger, N. Stoll and K. Thurow, "4D simulation system for laboratory workflow of life science automation," 2012 IEEE International Instrumentation and Measurement Technology Conference Proceedings, Graz, Austria, 2012, pp. 1886-1890, doi: 10.1109/I2MTC.2012.6229314.
- [10] S. Yu, L. Liu and M. Fu, "The Application Research on Knowledge Management of Project Manager," 2009 International Conference on Information Management, Innovation Management and Industrial Engineering, Xi'an, China, 2009, pp. 340-343, doi: 10.1109/ICIII.2009.391.
- [11] İ. Ç. Taş and A. A. Müngen, "Using Pre-Processing Methods to Improve OCR Performances of Digital Historical Documents," 2021 Innovations in Intelligent Systems and Applications Conference (ASYU), Elazig, Turkey, 2021, pp. 1-5, doi: 10.1109/ASYU52992.2021.9598972.
- [12] G. Cavalcanti and E. Filho, "An architecture for document management," Proceedings. International Conference on Image Processing, Rochester, NY, USA, 2002, pp. 973-976 vol.3, doi: 10.1109/ICIP.2002.1039137.
- [13] M. Ramirez, E. Tapia, M. Block and R. Rojas, "Quantile Linear Algorithm for Robust Binarization of Digitalized Letters," Ninth International Conference on Document Analysis and Recognition (ICDAR 2007), Curitiba, Brazil, 2007, pp. 1158-1162, doi: 10.1109/ICDAR.2007.4377097.
- [14] M. K. Ugale, S. J. Patil and V. B. Musande, "Document management system: A notion towards paperless office," 2017 1st International Conference on Intelligent Systems and Information Management (ICISIM), Aurangabad, India, 2017, pp. 217-224, doi: 10.1109/ICISIM.2017.8122176.
- [15] L. Schomaker, M. Bulacu and K. Franke, "Automatic writer identification using fragmented connected-component contours," Ninth International Workshop on Frontiers in Handwriting Recognition, Kokubunji, Japan, 2004, pp. 185-190, doi: 10.1109/IWFHR.2004.22.
- [16] D. Singh, S. V. S and V. K, "A Proposed Approach for Identifying the Connotative Relationship of English Sentences and paragraphs using the NLP package of Python," 2022 International Conference on Innovative Computing, Intelligent Communication and Smart Electrical Systems (ICES), Chennai, India, 2022, pp. 1-7, doi: 10.1109/ICES55317.2022.9914220.
- [17] M. Kalaiselvan and A. V. Kathiravan, "A pioneering tool for text summarization using star map," 2013 International Conference on Pattern Recognition, Informatics and Mobile Engineering, Salem, India, 2013, pp. 277-281, doi: 10.1109/ICPRIME.2013.6496486.
- [18] IBM. (2021, January 11). Optical Character Recognition (OCR): A Complete Guide. IBM Cloud Blog. [Online] Available:<https://www.ibm.com/cloud/blog/optical-character-recognition>. [Accessed: 7-Mar-2023].
- [19] AI Multiple. (n.d.). Handwriting Recognition: A Comprehensive Guide. [Online] Available:

- <https://research.aimultiple.com/handwriting-recognition/>. [Accessed: 3-Mar-2023].
- [20] Kundu, S. (2021). Handwritten Text Recognition: A Brief Review. In B. Gupta, D. Mandal, R. Kumar, & H. Gupta (Eds.), *Emerging Technologies for Agriculture and Environment* (pp. 309-318). Springer.
- [21] V. K. Vaisakh and L. B. Das, "Handwritten Malayalam Character Recognition System using Artificial Neural Networks," 2020 IEEE International Students' Conference on Electrical, Electronics and Computer Science (SCEECS), Bhopal, India, 2020, pp. 1-4, doi: 10.1109/SCEECS48394.2020.101.
- [22] M. Grüber, "Acoustic analysis of czech expressive recordings from a single speaker in terms of various communicative functions," 2011 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), Bilbao, Spain, 2011, pp. 293-298, doi: 10.1109/ISSPIT.2011.6151576.
- [23] S. M. Kim, C. J. Chun and H. K. Kim, "Multi-channel audio recording based on super directive beam forming for portable multimedia recording devices," in *IEEE Transactions on Consumer Electronics*, vol. 60, no. 3, pp. 429-435, Aug. 2014, doi: 10.1109/TCE.2014.6937327.
- [24] S. Das et al., "Hand-Written and Machine-Printed Text Classification in Architecture, Engineering & Construction Documents," 2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR), Niagara Falls, NY, USA, 2018, pp. 546-551, doi: 10.1109/ICFHR-2018.2018.00101.
- [25] F. Bourgault, T. Furukawa and H. F. Durrant-Whyte, "Process model, constraints, and the coordinated search strategy," *IEEE International Conference on Robotics and Automation*, 2004. *Proceedings. ICRA '04.* 2004, New Orleans, LA, USA, 2004, pp. 5256-5261, doi: 10.1109/ROBOT.2004.1302552.
- [26] F. M. Hanis, H. Khoshvaghti, M. Teimouri and H. Veisi, "A Language-Independent Approach to Classification of Textual File Fragments: Case Study of Persian, English, and Chinese Languages," 2021 11th International Conference on Computer Engineering and Knowledge (ICCKE), Mashhad, Iran, Islamic Republic of, 2021, pp. 254-259, doi: 10.1109/ICCKE54056.2021.9721512.
- [27] J. Sanguinetti, "High Level Design: The Future is Now," 2005 18th Symposium on Integrated Circuits and Systems Design, Florianopolis, Brazil, 2005, pp. 5-5, doi: 10.1109/SBCCI.2005.4286820.
- [28] G. Vamvakas, B. Gatos and S. J. Perantoni, "A Novel Feature Extraction and Classification Methodology for the Recognition of Historical Documents," 2009 10th International Conference on Document Analysis and Recognition, Barcelona, Spain, 2009, pp. 491-495, doi: 10.1109/ICDAR.2009.223.
- [29] A.Khalid, W. Q. Khan, R. Q. Khan and H. M. P. Memon, "Google image suggestions, extension for unfamiliar terminologies," 2017 International Conference on Information and Communication Technologies (ICICT), Karachi, Pakistan, 2017, pp. 16-22, doi: 10.1109/ICICT.2017.8320158.

Citation of this Article:

Thilakarathne U.T, Rajarathne R.M.H.M, Weerasinghe V.P, Kotuwegoda K.S.D, N.H.P. Ravi Supunya Swarnakantha, U.U. Samantha Rajapaksha, "Image Processing and Natural Language Processing Based Digitalized Document Generator" Published in *International Research Journal of Innovations in Engineering and Technology - IRJIET*, Volume 7, Issue 6, pp 123-130, June 2023. Article DOI <https://doi.org/10.47001/IRJIET/2023.706019>
